

Title	マルコフ型決定過程 (1) (動的計画法研究会報告集)
Author(s)	坂本, 武司
Citation	数理解析研究所講究録 (1967), 28: 75-85
Issue Date	1967-09
URL	<a href="http://hdl.handle.net/2433/107531">http://hdl.handle.net/2433/107531</a>
Right	
Type	Departmental Bulletin Paper
Textversion	publisher

## マルコフ型決定過程(I)

日大 生産工 坂本 武司

## §1. 平均最適政策

有限個の状態・決定からなる離散型マルコフ型決定過程について述べる。

$S$  個の状態  $1, 2, \dots, S$  からなる系 (System) を考える。離散時刻 ( $n = 1, 2, \dots$ ) で系の状態を観測し、可能な決定の集合  $K = \{k\} = \{1, 2, \dots, K\}$  から決定を選ぶ。もし、 $n$  時刻で状態  $i$  を観測し、決定  $k$  を行えば、その結果

(1) 系は推移確率  $q_{ij}^k$  によって新しい状態  $j$  に移る。

(2) 状態  $i$  で決定  $k$  を行えば、利得  $r_i^k$  を得る。

即ち、

状態空間:  $S = \{1, 2, \dots, S\}$

決定空間:  $K = \{1, 2, \dots, K\}$

推移確率:  $q_{ij}^k$  ,  $\sum_{j=1}^S q_{ij}^k = 1$

利得:  $r_i^k$

$f$  を  $S$  から  $K$  への関数とし、 $f$  の全体を  $F$  で表わす。 $F$  の要素の系列  $\pi = (f_1, f_2, \dots, f_n, \dots)$  を政策と呼ぶ。政策  $\pi$  を用いるとは初期状態が  $i$  であるば、 $\pi$  一期で決定  $f_1(i)$ 、 $\pi$  二期の状態が  $i_2$  であるば、 $f_2(i_2)$  以下同様にして、 $\pi$   $n$  期で  $f_n(i_n)$  を行う。 $\pi = (f, f, \dots, f, \dots) \equiv f^\infty$  を定常政策という。

$f \in F$  に対する利得及び推移確率を次のようにベクトル表示する。

$$r(f) = (r_1^{f(i)}, r_2^{f(i)}, \dots, r_s^{f(i)})' \quad Q(f) = (q_{ij}^{f(i)})$$

( $S \times 1$  ベクトル)                      ( $S \times S$  行列)

$V_n(i, \pi) =$  状態  $i$  から出発して、政策  $\pi$  を用いるときの  $n$  期迄の総期待利得

とする。そのベクトル表示を

$$V_n(\pi) = (V_n(1, \pi), \dots, V_n(s, \pi))' \quad (S \times 1 \text{ ベクトル})$$

とする。そのとき、

$$(1.1) \quad V_n(\pi) = \sum_{k=0}^{n-1} Q_k(\pi) r(f_{n+1})$$

但し、 $Q_k(\pi) = Q(f_1) \cdots Q(f_k)$ ,  $Q_0(\pi) = I$  (単位行列)

任意の  $S \times 1$  ベクトル  $W$  に対して作用素  $L(f)$  を

$$(1.2) \quad L(f)W = r(f) + Q(f)W$$

で定義する。 $L(f)$  は monotone である。即ち、任意のベクトル  $u, w$  について、すべての  $i$  について  $u_i \geq w_i$  ならば  $u \geq w$ 、 $u \geq w$  且つ  $u \neq w$  ならば  $u > w$  と定義すると、 $u \geq w$  且

らば、 $L(f)u \geq L(f)w$  である。今後上記のようにベクトルの順序を定めておく。 $L(f)$ を用いると、

$$V_n(\pi) = L(f_1)L(f_2)\cdots L(f_{n-1})r(f_n)$$

$V_n(\pi)$ は  $n \rightarrow \infty$  のとき発散し、1期当りの平均利得  $\frac{1}{n} V_n(\pi)$  は一般に収束しないが、 $u(\pi) = \liminf_{n \rightarrow \infty} \frac{1}{n} V_n(\pi)$  を考え、 $u(f) \geq u(\pi)$  for all  $\pi$  を満たす定常政策  $f^\infty$  が存在する。但し  $u(f) = u(f^\infty) = Q^*(f)r(f)$ 。また割引率  $\beta$  を入れると、 $\beta \rightarrow 1$  のとき、無限過程の総期待利得を最大にする定常政策が存在する。更に、各  $n$  について、 $U(n) = \max_{\pi} V_n(\pi)$  とすると、 $\lim_{n \rightarrow \infty} \frac{1}{n} U(n) = u(f)$  とする定常政策  $f^\infty$  が存在する。そして、この3つの定常政策の平均利得は一致する、そこで、定常政策の中で、平均利得を最大にする政策を決定する algorithm を  $\beta$ -変換及び割引率  $\beta$  を用いずには求めてみる。

Theorem 1. (1), (2). 任意の  $f \in F$  について、 $Q^n(f) \rightarrow Q^*(f)$  ( $n \rightarrow \infty$ ) とする。

$$(1.2) \quad V_n(f^\infty) = nu(f) + v(f) + \varepsilon(n, f)$$

(b).  $Q(f)$  のエルゴディク集合の周期の最小公倍数を  $N$  とし、

$$Q_0 \equiv Q^N \text{ とすると、}$$

$$(1.3) \quad V_{nN+m}(f^\infty) = (nN+m)u(f) + v(f) + w(m, f) + \varepsilon(n, f)$$

但し、 $u(f)$  は

$$(1.4) \quad (I - Q(f))u = 0, \quad Q^*(f)u = Q^*(f)r(f)$$

の一意の解であり、 $v(f)$ は

$$(1.5) \quad (I - Q(f))v = r(f) - u(f), \quad Q^*(f)v = 0$$

の一意の解である。また、 $w(nN+m, f) = w(m, f)$  ( $n=1, 2, \dots$ )

$$\varepsilon(n, f) \rightarrow 0 \quad (n \rightarrow \infty)$$

$$(2) \quad G_1(f) \equiv \{g \mid g \in F, Q(g)u(f) > u(f)\}$$

$$G_2(f) \equiv \{g \mid g \in F, Q(g)u(f) = u(f), r(g) + Q(g)v(f) > u(f) + v(f)\}$$

$$G(f) \equiv G_1(f) \cup G_2(f) \text{ とする.}$$

$$g \in G(f) \text{ ならば, } u(g) \geq u(f)$$

$$(3) \quad G(f) = \emptyset \text{ ならば, 任意の } g \in F \text{ について } u(f) \geq u(g)$$

$$(\text{証明}) \quad (1) \quad (2). \quad V_n(f^\infty) = \sum_{k=0}^{n-1} Q^k(f) r(f)$$

$$= (n-1) Q^*(f) r(f) + \sum_{k=0}^{n-1} (Q(f) - Q^*(f))^k r(f)$$

$n \rightarrow \infty$  のとき、 $Q^n(f) \rightarrow Q^*(f)$  となる。  $n$  が十分大になるとき、

$$V_n(f^\infty) = n Q^*(f) r(f) + H(f) r(f) + \varepsilon(n, f)$$

$$= n u(f) + v(f) + \varepsilon(n, f)$$

$$\text{但し, } v(f) \equiv H(f) r(f), \quad H(f) \equiv (I - Q(f) + Q^*(f))^{-1} - Q^*(f)$$

$$\varepsilon(n, f) \rightarrow 0 \quad (n \rightarrow \infty)$$

$$(b) \quad Q^*(f) = \frac{1}{N} Q_0^*(f) \sum_{i=0}^{N-1} Q^i(f) \text{ より}$$

$$V_{nN+m}(f^\infty) = (n-1) N Q^*(f) r(f) + [I - (Q_0(f) - Q_0^*(f))]^{-1} \sum_{i=0}^{N-1} Q^i(f) r(f) + Q_0^*(f) \sum_{i=0}^{m-1} Q^i(f) r(f)$$

よって、2 より上の結果を得る。

(2). 周期的な場合には  $\alpha = 2$  示す. 仮定より  $n$  を十分大きくとると.

$$\begin{aligned} V_{nN+m+1}(g, f^\infty) &= Q(g)W(m, f) - Q(g)\varepsilon(nN+m, f) \\ &> V_{nN+m+1}(f^\infty) - W(m+1, f) - \varepsilon(nN+m+1, f) \end{aligned}$$

即ち. 
$$L(g)V_{nN+m}(f^\infty) > V_{nN+m+1}(f^\infty) + Q(g)W(m, f) - \varepsilon(nN+m+1, f) + Q(g)\varepsilon(nN+m, f)$$

一般に、任意の自然数  $M$  について.

$$\begin{aligned} L^M(g)V_{nN+m}(f^\infty) &> V_{nN+m+M}(f^\infty) + Q^M(g)W(m, f) \\ &\quad - W(m+M, f) - \varepsilon(nN+m+M, f) + Q^M(g)\varepsilon(nN+m, f) \end{aligned}$$

実際.  $M=1$  のとき明らかに  $M$  のとき仮定すると.

$$\begin{aligned} L^{M+1}(g)V_{nN+m}(f^\infty) &\geq L(g)V_{nN+m+M}(f^\infty) + Q^{M+1}(g)W(m, f) \\ &\quad - Q(g)W(m+M, f) - Q(g)\varepsilon(nN+m+M, f) \\ &\quad + Q^{M+1}(g)\varepsilon(nN+m, f) \end{aligned}$$

仮定を用いると.

$$\begin{aligned} L^{M+1}(g)V_{nN+m}(f^\infty) &> V_{nN+m+1}(f^\infty) + Q^{M+1}(g)W(m, f) \\ &\quad - W(m+M+1, f) - \varepsilon(nN+m+M+1, f) + Q^{M+1}(g)\varepsilon(nN+m, f) \end{aligned}$$

$$\min_{\substack{1 \leq i \leq S \\ 1 \leq m \leq N}} W_i(m, f) = C_1, \quad \max_{\substack{1 \leq i \leq S \\ 1 \leq m \leq N}} W_i(m, f) = C_2, \quad \delta = (1, 1, \dots, 1)'$$

とおく.

$$\begin{aligned} L^M(g)V_{nN+m}(f^\infty) &> V_{nN+m+M}(f^\infty) + (C_1 - C_2)\delta \\ &\quad - \varepsilon(nN+m+M, f) + Q^M(g)\varepsilon(nN+m, f) \end{aligned}$$

$M \rightarrow \infty$  とすると.  $\varepsilon(nN+m+M, f) \rightarrow 0$  かつ  $\varepsilon(nN+m, f)$

$Q^M$  は有界だから、上式を  $M$  で割り  $M \rightarrow \infty$  とすると、 $u(f) \geq u(f)$  を得る。同期的でないときも同様である。

## 8.2. 強最適政策の計算法

割引率  $\beta$  ( $0 \leq \beta < 1$ ) を考えたときの政策  $\pi$  に対する終期待利得を  $V_\beta(\pi)$  で表わすと、

$$(2.1) \quad V_\beta(\pi) = \sum_{n=0}^{\infty} \beta^n Q_n(\pi) r(f_{n+1}) \quad \text{但し} \quad Q_0(\pi) = I$$

もし、 $U(\beta) \equiv V_\beta(\pi^*) \geq V_\beta(\pi)$  for all  $\pi$  が成立するとき、 $\pi^*$  を  $\beta$ -最適政策と呼ぶ。

また、適当な  $\beta_0$  が存在して ( $0 \leq \beta_0 < 1$ )、すべて  $\beta_0 \leq \beta < 1$  なる  $\beta$  に対して、 $V_\beta(\pi^*) \geq V_\beta(\pi)$  for all  $\pi$  が成立するとき、 $\pi^*$  を最適政策と呼ぶ。 $\beta$ -最適及び最適な定常政策が存在するとは知られている。

また、 $U(\beta) - V_\beta(\pi) \rightarrow 0$  ( $\beta \uparrow 1$ ) ならば  $\pi$  を強最適 (nearly optimal) 又は 1-最適 (1-optimal) と呼ぶ。2.2 では、強最適政策を決定する Veinott の algorithm について述べる。

まず、準備として D. Blackwell の結果を述べる。

Theorem 1 と同じ手順により

Lemma 1. 任意の  $f \in F$  について

$$(2.2) \quad V_\beta(f^\infty) = \frac{u(f)}{1-\beta} + v(f) + \varepsilon(\beta, f), \quad 0 \leq \beta < 1$$

$$\varepsilon(\beta, f) \rightarrow 0 \quad (\beta \rightarrow 1)$$

次の集合を定義する。

$$F' \equiv \{ f \mid f \in F, u(f) \geq u(g) \text{ all } g \in F \}$$

$$F'' \equiv \{ f \mid f \in F', v(f) \geq v(g) \text{ all } g \in F' \}$$

(2.2)式の形より

Lemma 2.  $F''$  は殆ど最適なすべての  $f \in F$  の集合である。

したがって  $F''$  の要素を決定する algorithm を作らばよい。  
定理 1 より  $G(f) = \emptyset$  ならば  $f \in F'$  であるから、 $F'$  の中で  $v$  を最大にする計算法を求めよばよい。

Lemma 3.  $f \in F, g \in G(f)$  ならば:  $u(g) > u(f)$  又は、

$$u(g) = u(f), v(g) > v(f)$$

(証明) 定理 1 より、 $g \in G(f)$  ならば  $u(g) \geq u(f)$ 。又:  $u(g) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n Q^k(f) r(f) = \lim_{\beta \rightarrow 1} (1-\beta) \sum_{k=0}^{\infty} \beta^k Q^k(f) r(f) = \lim_{\beta \rightarrow 1} (1-\beta) V_{\beta}(f^{\infty})$  であり、十分大に近い  $\beta$  に対して:  $V_{\beta}(g^{\infty}) > V_{\beta}(f^{\infty})$

$$V_{\beta}(g^{\infty}) = \frac{u(g)}{1-\beta} + v(g) + \varepsilon(\beta, g)$$

$$V_{\beta}(f^{\infty}) = \frac{u(f)}{1-\beta} + v(f) + \varepsilon(\beta, f)$$

であるから、 $u(g) \geq u(f)$  且つ、 $u(g) = u(f)$  ならば  $v(g) \geq v(f)$

もし、 $(u(g), v(g)) = (u(f), v(f))$  ならば、 $g \in G(f)$  は矛盾。

次の集合を定義する。

$$E(f) = \{ g \mid g \in F, Q(g)u(f) = u(f), r(g) + Q(g)v(f) = u(f) + v(f) \}$$

$E(f)$  の定義と (1.3), (1.4) より  $g \in E(f)$  ならば  $u(g) = u(f) \in T$  である。



1.  $f \in F$  ならば  $E(f) \subset F$  である。集合  $E(f)$  は計算結果から  $E(f)$  の中で  $v$  を最大にする  $z$  を与える。

Lemma 4. (Veinott)  $f \in F, g \in E(f)$  とする。  $w(g)$  は

$$(2.3) \quad [I - Q(g)] w = 0, \quad Q^*(g) w = Q^*(g) (-v(f))$$

の一意の解とする。そのとき

$$(2.4) \quad v(g) = v(f) + w(g)$$

(証明) 一意性は定理 2 より出る。  $g \in E(f)$  より

$$[I - Q(g)] v(f) = r(g) - u(f)$$

この式と (2.3) を加えると

$$[I - Q(g)] [v(f) + w(g)] = r(g) - u(f), \quad Q^*(g) [v(f) + w(g)] = 0$$

よって定理 2 より、(2.4) が出る。

この Lemma より  $E(f)$  の中で  $w(g)$  を最大にする  $z$  とは、 $F$  を  $E(f)$  で  $r(f)$  を  $-v(f)$  でおきかえれば、 $F$  の中で  $u(f)$  を最大にするのと同じ形をしていふ  $z$  とわかる。従って  $E(f)$  の中で  $v(g)$  を最大にするためには、政策反復法を再び使えばよいことになる。

任意の  $f \in F$  について、  $z(f)$  を (2.5) の一意の解とする。

$$(2.5) \quad (I - Q(f)) z(f) = -v(f), \quad Q^*(f) z(f) = 0$$

$$(2.6) \quad H(f) \equiv \{ g \mid g \in E(f), -v(f) + Q(g) z(f) > z(f) \}$$

とする。そのとき

Theorem 2. (Veinott)

(1)  $f \in F$ ,  $G(f) = \emptyset$  とする。もし、 $v(f) \geq v(g)$  for all

$g \in E(f)$  ならば  $f \in F''$

(2)  $G(f) \cup H(f) = \emptyset$  ならば  $f \in F''$

(3)  $g \in H(f)$  ならば  $u(g) = u(f)$  であり、更に  $v(g) > v(f)$

又は  $v(g) = v(f)$ ,  $z(g) > z(f)$  のいずれか成立する。

(証明) (1)  $g \in F' - E(f)$  とする。  $G(f) = \emptyset$ ,  $g \in F'$  より  $u(g)$

$= u(f)$ . 更に  $g \notin E(f)$  より  $r(g) + Q(g)v(f) < u(f) + v(f)$

$$V_\beta(g, f^\infty) = \frac{Q(g)u(f)}{1-\beta} + r(g) - Q(g)u(f) + Q(g)v(f) + \varepsilon(\beta, f, g)$$

$$V_\beta(f^\infty) = \frac{u(f)}{1-\beta} + v(f) + \varepsilon(\beta, f)$$

であるから、十分小さい  $\beta$  に対して  $V_\beta(g, f^\infty) < V_\beta(f^\infty)$

より  $V_\beta(g^\infty) < V_\beta(f^\infty)$  ([1]参照). 即ち  $v(f) > v(g)$ .

(2)  $H(f) = \emptyset$  であるから、 $w(f) \geq w(g)$  for all  $g \in E(f)$ . 然

るに  $w(f) = 0$  であるから、 $v(g) = v(f) + w(g) \leq v(f) + w(f) = v(f)$  for all  $g \in E(f)$ .

(3)  $g \in H(f) \subset E(f)$  より  $u(g) = u(f)$ . 更に Lemma 3 より、

$w(g) > w(f) = 0$ . 又は  $w(g) = w(f) = 0$ ,  $z(g) > z(f)$ .

この定理により、 $f_1 \in F$  とし、 $f_2, f_3, \dots$  を  $f_{i+1} \in G(f_i) \cup H(f_i)$  であるように選ぶ。  $\{u_i(f), v_i(f), z_i(f)\}$  は辞書的順序で増加するから、同じ  $f_i$  が2度現われることはない。  $F$  は有限集合だから、ある  $i$  で  $G(f_i) \cup H(f_i) = \emptyset$  となり、殆ど最適政策が得られる。

例.  $S = \{1, 2\}$ ,  $K_1 = \{1, 2\}$ ,  $K_2 = \{1\}$  ( $K_i$  は状態  $i$  2-可能な決定の場合)

$$F = \{f, g\} \text{ 且 } f(i) = 1 \quad (i=1, 2)$$

$$g(1) = 2, \quad g(2) = 1$$

$$r(f) = \begin{pmatrix} 3 \\ -3 \end{pmatrix} \quad Q(f) = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad Q^*(f) = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

$$r(g) = \begin{pmatrix} 6 \\ -3 \end{pmatrix} \quad Q(g) = \begin{pmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad Q^*(g) = \begin{pmatrix} \frac{1}{3} & \frac{2}{3} \\ \frac{1}{3} & \frac{2}{3} \end{pmatrix}$$

$$u(f) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad v(f) = \begin{pmatrix} 3 \\ -3 \end{pmatrix} \quad z(f) = \begin{pmatrix} -3 \\ 3 \end{pmatrix}$$

$$u(g) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad v(g) = \begin{pmatrix} 4 \\ -2 \end{pmatrix} \quad z(g) = \begin{pmatrix} -\frac{8}{3} \\ \frac{4}{3} \end{pmatrix}$$

$$G(f) = \phi, \quad G(g) = \phi, \quad H(f) = g, \quad H(g) = \phi$$

$g$  が殆ど最適政策

### 参考文献

(1) D. Blackwell, Discrete Dynamic Programming, Ann. Math. Statist. 33 (1962) 719-726

(2) B. W. Brown, On the Iterative Method of Dynamic Programming on a Finite Space Discrete Time Markov Processes, Ann. Math. Statist. 36 (1965) 1279-1285

(3) C. Derman, On Sequential Decisions and Markov

- chains, Management Science 9 (1962) 16-24
- (4) R.A. Howard, Dynamic Programming and Markov Processes, 1960, Technology Press and Wiley
- (5) M. Ogawara, A Note on Discrete Markovian Decision Processes, Bull. Math. Statist 11 (1963) 35-42
- (6) A.F. Veinott, Jr, On Finding Optimal Policies in Discrete Dynamic Programming, Ann. Math. Statist 37 (1966) 1284-1294

